

107563065

IAP20 Rec'd PCT/PTO 30 DEC 2005

CERTIFICATION OF TRANSLATION

I, Bertrand LOISEL, of CABINET PLASSERAUD, 65/67 rue de la Victoire, 75440 PARIS CEDEX 09, FRANCE, do hereby declare that I am well acquainted with the English language, and attest that the document attached is a true English language translation of the text of International Patent Application no.PCT/FR03/02037.

Dated this 15<sup>th</sup> of December, 2005

Bertrand LOISEL

A handwritten signature in black ink, consisting of a large, stylized 'B' followed by a series of loops and a final vertical stroke.

IAP20 Rec'd PCT/PTO 30 DEC 2005

**Method and system for analysis of vocal signals for a  
compressed representation of speakers**

The present invention relates to a method and a device for analyzing vocal signals.

- 5 The analysis of vocal signals requires in particular the ability to represent a speaker. The representation of a speaker by a mixture of Gaussians ("Gaussian Mixture Model" or GMM) is an effective representation of the acoustic or vocal identity of a speaker.
- 10 According to this technique, the speaker is represented, in an acoustic reference space of a predetermined dimension, by a weighted sum of a predetermined number of Gaussians.

- This type of representation is accurate when a large
- 15 amount of data is available, and when there are no physical constraints in respect of the storage of the parameters of the model, or in respect of the execution of the calculations on these numerous parameters.

- Now, in practice, to represent a speaker within IT
- 20 systems, it transpires that the time for which a speaker is talking is short, and that the size of the memory required for these representations, as well as the times for calculations with regard to these parameters are too big.

- 25 It is therefore important to seek to represent a speaker in such a way as to drastically reduce the number of parameters required for the representation thereof while maintaining correct performance. Performance is meant as the error rate of vocal
- 30 sequences that are not recognized as belonging or not to a speaker with respect to the total number of vocal sequences.

Solutions in this regard have been proposed, in particular in the document "SPEAKER INDEXING IN LARGE

AUDIO DATABASES USING ANCHOR MODELS" by D.E. Sturim, D.A. Reynolds, E. Singer and J.P. Campbell. Specifically, the authors propose that a speaker be represented not in an absolute manner in an acoustic reference space, but instead in a relative manner with respect to a predetermined set of representations of reference speakers also called anchor models, for which GMM-UBM models are available (UBM standing for "Universal Background Model"). The proximity between a speaker and the reference speakers is evaluated by means of a Euclidean distance. This enormously decreases the calculational load, but the performance is still limited and inadequate.

In view of the foregoing, an object of the invention is to analyze vocal signals by representing the speakers with respect to a predetermined set of reference speakers, with a reduced number of parameters decreasing the calculational load for real-time applications, with acceptable performance, by comparison with analysis using a representation by the GMM-UBM model.

It is then for example possible to perform indexings of audio documents of large databases where the speaker is the indexing key.

Thus, according to an aspect of the invention, there is proposed a method of analyzing vocal signals of a speaker ( $\lambda$ ), using a probability density representing the resemblances between a vocal representation of the speaker ( $\lambda$ ) in a predetermined model and a predetermined set of vocal representations of a number  $E$  of reference speakers in said predetermined model, and the probability density is analyzed so as to deduce therefrom information on the vocal signals.

This makes it possible to drastically decrease the number of parameters used, and allows devices

implementing this method to be able to work in real time, while decreasing the calculation time, while decreasing the size of the memory required.

In a preferred embodiment, an absolute model (GMM), of dimension D, using a mixture of M Gaussians, is taken as predetermined model, for which the speaker ( $\lambda$ ) is represented by a set of parameters comprising weighting coefficients ( $\alpha_i$ ,  $i = 1$  to M) for the mixture of Gaussians in said absolute model (GMM), mean vectors ( $\mu_i$ ,  $i = 1$  to M) of dimension D and covariance matrices ( $\Sigma_i$ ,  $i = 1$  to M) of dimension  $D \times D$ .

In an advantageous embodiment, the probability density of the resemblances between the representation of said vocal signals of the speaker ( $\lambda$ ) and the predetermined set of vocal representations of the reference speakers is represented by a Gaussian distribution ( $\psi(\mu^\lambda, \Sigma^\lambda)$ ) of mean vector ( $\mu^\lambda$ ) of dimension E and of covariance matrix ( $\Sigma^\lambda$ ) of dimension  $E \times E$  which are estimated in the space of resemblances to the predetermined set of E reference speakers.

In a preferred embodiment, the resemblance ( $\psi(\mu^\lambda, \Sigma^\lambda)$ ) of the speaker ( $\lambda$ ) with respect to the E reference speakers is defined, for which speaker ( $\lambda$ ) there are  $N_\lambda$  segments of vocal signals represented by  $N_\lambda$  vectors of the space of resemblances with respect to the predetermined set of E reference speakers, as a function of a mean vector ( $\mu^\lambda$ ) of dimension E and of a covariance matrix ( $\Sigma^\lambda$ ) of the resemblances of the speaker ( $\lambda$ ) with respect to the E reference speakers.

In an advantageous embodiment, a priori information is moreover introduced into the probability densities of the resemblances ( $\psi(\tilde{\mu}^\lambda, \tilde{\Sigma}^\lambda)$ ) with respect to the E reference speakers.

In a preferred embodiment, the covariance matrix of the

speaker ( $\lambda$ ) is independent of said speaker ( $\tilde{\Sigma}^\lambda = \tilde{\Sigma}$ ).

According to another aspect of the invention, there is proposed a system for the analysis of vocal signals of a speaker ( $\lambda$ ), comprising databases in which are stored  
5 vocal signals of a predetermined set of E reference speakers and their associated vocal representations in a predetermined model, as well as databases of audio archives, characterized in that it comprises means of analysis of the vocal signals using a vector  
10 representation of the resemblances between the vocal representation of the speaker and the predetermined set of vocal representations of E reference speakers.

In an advantageous embodiment, the databases also store the vocal signals analysis performed by said means of  
15 analysis.

The invention may be applied to the indexing of audio documents, however other applications may also be envisaged, such as the acoustic identification of a speaker or the verification of the identity of a  
20 speaker.

Other objects, features and advantages of the invention will become apparent on reading the following description, given by way of nonlimiting example, and offered with reference to the single appended drawing  
25 illustrating an application of a use of the method in respect of the indexing of audio documents.

The figure represents an application of the system according to an aspect of the invention in respect of the indexing of audio databases. Of course, the  
30 invention applies also to the acoustic identification of a speaker or the verification of the identity of a speaker, that is to say, in a general manner, to the recognition of information relating to the speaker in the acoustic signal. The system comprises a means for

receiving vocal data of a speaker, for example a mike  
 1, linked by a wire or wireless connection 2 to means  
 of recording 3 of a request enunciated by a speaker  $\lambda$   
 and comprising a set of vocal signals. The recording  
 5 means 3 are linked by a connection 4 to storage means 5  
 and, by a connection 6, to means of acoustic processing  
 7 of the request. These acoustic means of processing  
 transform the vocal signals of the speaker  $\lambda$  into a  
 representation in an acoustic space of dimension  $D$  by a  
 10 GMM model for representing the speaker  $\lambda$ .

This representation is defined by a weighted sum of  $M$   
 Gaussians according to the equations:

$$p(x|\lambda) = \sum_{i=1}^M \alpha_i b_i(x) \quad (1)$$

$$b_i(x) = \frac{1}{(2\pi)^{D/2} |\Sigma_i|^{1/2}} \times \exp\left[-\frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i)\right] \quad (2)$$

$$\sum_{i=1}^M \alpha_i = 1 \quad (3)$$

in which:

15  $D$  is the dimension of the acoustic space of the  
 absolute GMM model;

$x$  is an acoustic vector of dimension  $D$ , i.e. vector of  
 the cepstral coefficients of a vocal signal sequence of  
 the speaker  $\lambda$  in the absolute GMM model;

20  $M$  denotes the number of Gaussians of the absolute GMM  
 model, generally a power of 2 lying between 16 and  
 1024;

$b_i(x)$  denotes, for  $i = 1$  to  $D$ , Gaussian densities,  
 parameterized by a mean vector  $\mu_i$  of dimension  $D$  and a  
 25 covariance matrix  $\Sigma_i$  of dimension  $D \times D$ ; and

$\alpha_i$  denotes, for  $i = 1$  to  $D$ , the weighting coefficients  
 of the mixture of Gaussians in the absolute GMM model.

The means of acoustic processing 7 of the request are

linked by a connection 8 to means of analysis 9. These means of analysis 9 are able to represent a speaker by a probability density vector representing the resemblances between the vocal representation of said speaker in the GMM model chosen and vocal representations of E reference speakers in the GMM model chosen. The means of analysis 9 are furthermore able to perform tests for verifying and/or identifying a speaker.

10 To carry out these tests, the analysis means undertake the formulation of the vector of probability densities, that is to say of resemblances between the speaker and the reference speakers.

This entails describing a relevant representation of a single segment  $x$  of the signal of the speaker  $\lambda$  by means of the following equations:

$$w^\lambda = \begin{pmatrix} \tilde{p}(x^\lambda | \bar{\lambda}_1) \\ \vdots \\ \tilde{p}(x^\lambda | \bar{\lambda}_E) \end{pmatrix} \quad (4)$$

$$\tilde{p}(x^\lambda | \bar{\lambda}_j) = \frac{1}{T_x} \log \left( \frac{p(x^\lambda | \bar{\lambda}_j)}{p(x^\lambda | \bar{\lambda}_{UBM})} \right) \quad (5)$$

$$p(x | \bar{\lambda}) = \sum_{k=1}^M \alpha_k b_k(x) \text{ where } \sum_{k=1}^M \alpha_k = 1 \quad (6)$$

$$b_k(x) = \frac{1}{(2\pi)^{D/2} |\Sigma_k|^{1/2}} \times \exp \left[ -\frac{1}{2} {}^t(x - \mu_k)(\Sigma_k)^{-1}(x - \mu_k) \right] \quad (7)$$

in which:

$w^\lambda$  is a vector of the space of resemblances to the predetermined set of E reference speakers representing the segment  $x$  in this representation space;

$\tilde{p}(x^\lambda | \bar{\lambda}_j)$  is a probability density or probability normalized by a universal model, representing the resemblance of the acoustic representation  $x^\lambda$  of a

segment of vocal signal of a speaker  $\lambda$ , given a reference speaker  $\bar{\lambda}_j$ ;

$T_x$  is the number of frames or of acoustic vectors of the speech segment  $x$ ;

- 5  $p(x^\lambda | \bar{\lambda}_j)$  is a probability representing the resemblance of the acoustic representation  $x^\lambda$  of a segment of vocal signal of a speaker  $\lambda$ , given a reference speaker  $\bar{\lambda}_j$ ;

- $p(x^\lambda | \bar{\lambda}_{UBM})$  is a probability representing the resemblance of the acoustic representation  $x^\lambda$  of a segment of vocal signal of a speaker  $\lambda$  in the model of the UBM world;
- 10

$M$  is the number of Gaussians of the relative GMM model, generally a power of 2 lying between 16 and 1024;

$D$  is the dimension of the acoustic space of the absolute GMM model;

- 15  $x^\lambda$  is an acoustic vector of dimension  $D$ , i.e. a vector of the cepstral coefficients of a sequence of vocal signal of the speaker  $\lambda$  in the absolute GMM model;

- $b_k(x)$  represents, for  $k = 1$  to  $D$ , Gaussian densities, parameterized by a mean vector  $\mu_k$  of dimension  $D$  and a covariance matrix  $\Sigma_k$  of dimension  $D \times D$ ;
- 20

$\alpha_k$  represents, for  $k = 1$  to  $D$ , the weighting coefficients of the mixture of Gaussians in the absolute GMM model.

- On the basis of the representations  $W_j$  of the segments of speech  $x_j$  ( $j = 1, \dots, N_\lambda$ ) of the speaker  $\lambda$ , the speaker  $\lambda$  is represented by the Gaussian distribution  $\psi$  of parameters  $\mu^\lambda$  and  $\Sigma_\lambda$  defined by the following relations:
- 25

$$\left\{ \begin{array}{l} \mu^\lambda = \{\mu_i^\lambda\}_{i=1, \dots, E} \quad \text{with} \quad \mu_i^\lambda = \frac{1}{N_\lambda} \sum_{j=1}^{N_\lambda} \tilde{p}(x_j^\lambda | \bar{\lambda}_i) \end{array} \right. \quad (8)$$

$$\left\{ \begin{array}{l} \Sigma^\lambda = \{\Sigma_{i'i'}^\lambda\}_{i,i'=1, \dots, E} \quad \text{with} \quad \Sigma_{i'i'}^\lambda = \frac{1}{N_\lambda} \sum_{j=1}^{N_\lambda} (\tilde{p}(x_j^\lambda | \bar{\lambda}_i) - \mu_i^\lambda)(\tilde{p}(x_j^\lambda | \bar{\lambda}_{i'}) - \mu_{i'}^\lambda) \end{array} \right. \quad (9)$$



in which  $\mu_i^\lambda$  represents components of the mean vector  $\mu^\lambda$  of dimension E of the resemblances  $\psi(\mu^\lambda, \Sigma^\lambda)$  of the speaker  $\lambda$  with respect to the E reference speakers, and  $\Sigma_{ii}^\lambda$  represents components of the covariance matrix  $\Sigma^\lambda$  of dimension E  $\times$  E of the resemblances  $\psi(\mu^\lambda, \Sigma^\lambda)$  of the speaker  $\lambda$  with respect to the E reference speakers.

The analysis means 9 are linked by a connection 10 to training means 11 making it possible to calculate the vocal representations, in the form of vectors of dimension D, of the E reference speakers in the GMM model chosen. The training means 11 are linked by a connection 12 to a database 13 comprising vocal signals of a predetermined set of speakers and their associated vocal representations in the reference GMM model. This database may also store the result of the analysis of vocal signals of initial speakers other than said E reference speakers. The database 13 is linked by the connection 14 to the means of analysis 9 and by a connection 15 to the acoustic processing means 7.

The system further comprises a database 16 linked by a connection 17 to the acoustic processing means 7, and by a connection 18 to the analysis means 9. The database 16 comprises audio archives in the form of vocal items, as well as the associated vocal representations in the GMM model chosen. The database 16 is also able to store the associated representations of the audio items calculated by the analysis means 9. The training means 11 are furthermore linked by a connection 19 to the acoustic processing means 7.

An example will now be described of the manner of operation of this system that can operate in real time since the number of parameters used is appreciably reduced with respect to the GMM model, and since many steps may be performed off-line.

The training module 11 will determine the

representations in the reference GMM model of the E reference speakers by means of the vocal signals of these E reference speakers stored in the database 13, and of the acoustic processing means 7. This  
5 determination is performed according to relations (1) to (3) mentioned above. This set of E reference speakers will represent the new acoustic representation space. These representations of the E reference speakers in the GMM model are stored in memory, for  
10 example in the database 13. All this may be performed off-line.

When vocal data are received from a speaker  $\lambda$ , for example via the mike 1, they are transmitted via the connection 2 to the recording means 3 able to perform  
15 the storage of these data in the storage means 5 with the aid of the connection 4. The recording means 3 transmit this recording to the means of acoustic processing 7 via the connection 6. The means of acoustic processing 7 calculate a vocal representation  
20 of the speaker in the predetermined GMM model as set forth earlier with reference to the above relations (1) to (3).

Furthermore, the means of acoustic processing 7 have calculated, for example off-line, the vocal  
25 representations of a set of S test speakers and of a set of T speakers in the predetermined GMM model. These sets are distinct. These representations are stored in the database 13. The means of analysis 9 calculate, for example off-line, a vocal representation of the S  
30 speakers and of the T speakers with respect to the E reference speakers. This representation is a vector representation with respect to these E reference speakers, as described earlier. The means of analysis 9 also perform, for example off-line, a vocal  
35 representation of the S speakers and of the T speakers with respect to the E reference speakers, and a vocal

representation of the items of the speakers of the audio base. This representation is a vector representation with respect to these E reference speakers.

5 The processing means 7 transmit the vocal representation of the speaker  $\lambda$  in the predetermined GMM model to the means of analysis 9, which calculate a vocal representation of the speaker  $\lambda$ . This representation is a representation by probability  
10 density of the resemblances to the E reference speakers. It is calculated by introducing a priori information by means of the vocal representations of T speakers. Specifically, the use of this a priori information makes it possible to maintain a reliable  
15 estimate, even when the number of available speech segments of the speaker  $\lambda$  is small. A priori information is introduced by means of the following equations:

$$\left\{ \begin{array}{l} \bar{\mu}^{\lambda} = \frac{N_0 \mu_0 + N_{\lambda} \mu^{\lambda}}{N_0 + N_{\lambda}} \end{array} \right. \quad (10)$$

$$\left\{ \begin{array}{l} W = \left( w_1^{\text{spk}-1} \dots w_{N_i}^{\text{spk}-1} \dots w_1^{\text{spk}-T} \dots w_{N_T}^{\text{spk}-T} \right) \end{array} \right. \quad (11)$$

20 in which:

$\mu^{\lambda}$ : mean vector of dimension E of the resemblances  
 $\psi(\mu^{\lambda}, \Sigma^{\lambda})$  of the speaker  $\lambda$  with respect to the E reference speakers;

$N_{\lambda}$ : number of segments of vocal signals of the speaker  
25  $\lambda$ , represented by  $N_{\lambda}$  vectors of the space of resemblances to the predetermined set of E reference speakers;

W: matrix of all the initial data of a set of T speakers  $\text{spk}_i$ , for  $i = 1$  to T, whose columns are  
30 vectors of dimension E representing a segment of

vocal signal represented by a vector of the space of resemblances to the predetermined set of E reference speakers, each speaker  $\text{spk}_i$  having  $N_i$  vocal segments, characterized by its mean vector  $\mu_0$  of dimension E, and by its covariance matrix  $\Sigma_0$  of dimension  $E \times E$ ;

$\bar{\mu}^\lambda$ : mean vector of dimension E of the resemblances  $\psi(\bar{\mu}^\lambda, \bar{\Sigma}^\lambda)$  of the speaker  $\lambda$  with respect to the E reference speakers, with introduction of a priori information; and

$\bar{\Sigma}^\lambda$ : covariance matrix of dimension  $E \times E$  of the resemblances  $\psi(\bar{\mu}^\lambda, \bar{\Sigma}^\lambda)$  of the speaker  $\lambda$  with respect to the E reference speakers with introduction of a priori information.

Moreover, it is possible to take a single covariance matrix for each speaker, thereby making it possible to orthogonalize said matrix off-line, and the calculations of probability densities will then be performed with diagonal covariance matrices. In this case, this single covariance matrix is defined according to the relations:

$$\begin{cases} \tilde{\Sigma}_{ij} = \frac{1}{N_0} \sum_{s=1}^T \sum_{j \in I_s} (W_{ij} - \bar{W}_{is})(W_{i'j} - \bar{W}_{i's}) & (12) \\ \bar{W}_{is} = \frac{1}{N_T} \sum_{j \in I_s} W_{ij} & (13) \end{cases}$$

in which

W is a matrix of all the initial data of a set of T speakers  $\text{spk}_i$ , for  $i = 1$  to T, whose columns are vectors of dimension E representing a segment of vocal signal represented by a vector of the space of resemblances to the predetermined set of E reference speakers, each speaker  $\text{spk}_i$  having  $N_i$  vocal segments, characterized by its mean vector  $\mu_0$  of dimension E, and by its covariance matrix  $\Sigma_0$  of dimension  $E \times E$ .

Next, the analysis means 9 will compare the vocal representations of the request and of the items of the base by identification and/or verification tests of the speakers. The speaker identification test consists in  
5 evaluating a measure of likelihood between the vector of the test segment  $w_x$  and the set of representations of the items of the audio base. The speaker identified corresponds to the one which gives a maximum likelihood score, i.e.  $\hat{\lambda} = \arg \max_{\lambda} p(w_x | \tilde{\mu}^{\lambda}, \tilde{\Sigma}^{\lambda})$  (14) from among the  
10 set of S speakers.

The speaker verification test consists in calculating a score of likelihood between the vector of the test segment  $w_x$  and the set of representations of the items of the audio base, normalized by its score of  
15 likelihood with the representation of the a priori information. The segment is authenticated if the score exceeds a predetermined given threshold, said score being given by the following relation:

$$\text{score} = \frac{p(w_x | \tilde{\mu}^{\lambda}, \tilde{\Sigma}^{\lambda})}{p(w_x | \mu_0, \Sigma_0)} \quad (15)$$

20 Each time the speaker  $\lambda$  is recognized in an item of the base, this item is indexed by means of information making it possible to ascertain that the speaker  $\lambda$  is talking in this audio item.

This invention can also be applied to other uses, such  
25 as the recognition or the identification of a speaker.

This compact representation of a speaker makes it possible to drastically reduce the calculation cost, since there are many fewer elementary operations in view of the drastic reduction in the number of  
30 parameters required for the representation of a speaker.

For example, for a request of 4 seconds of speech of a speaker, that is to say 250 frames, for a GMM model of dimension 27, with 16 Gaussians the number of elementary operations is reduced by a factor of 540, 5 thereby enormously reducing the calculation time. Furthermore, the size of memory used to store the representations of the speakers is appreciably reduced.

The invention therefore makes it possible to analyze vocal signals of a speaker while drastically reducing 10 the time for calculation and the memory size for storing the vocal representations of the speakers.